



# TENSOR CORE PROGRAMMABILITY AND PROFILING FOR AI AND HPC APPLICATIONS

Griffin Lacey

Max Katz

# TOO LONG; DIDN'T LISTEN

- Tensor Cores enable fast **mixed precision matrix multiplications**
- Growing number of AI/HPC examples **accelerated up to 25x**
- Mature software support with **high-level APIs and Nsight developer tools**
- All you need is **Volta / Turing GPU**

# OUTLINE

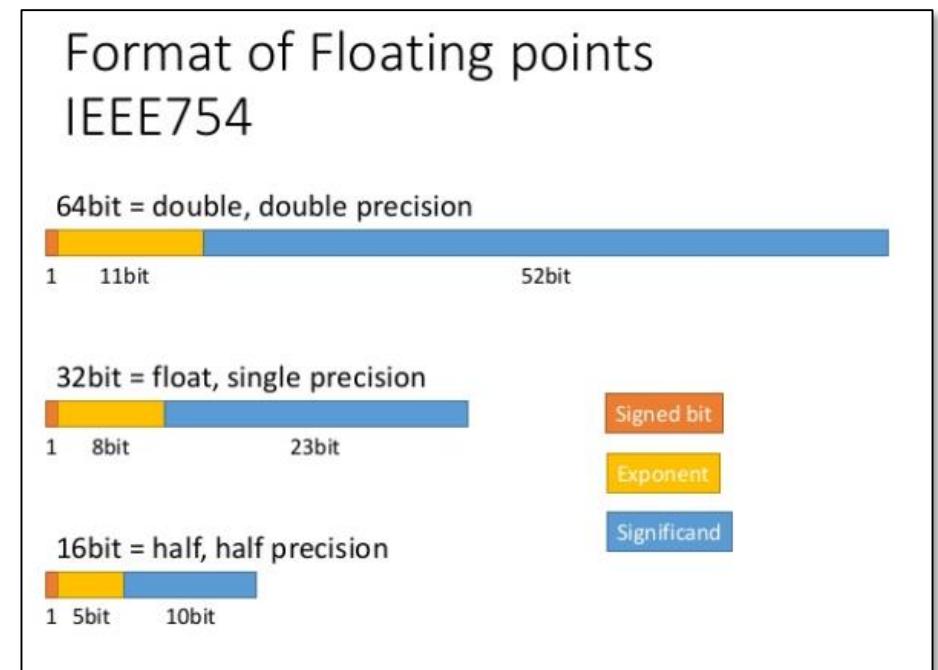
1. What are Tensor Cores?
2. Tensor Cores for AI
3. Tensor Cores for HPC
4. Profiling Tensor Cores

# OUTLINE

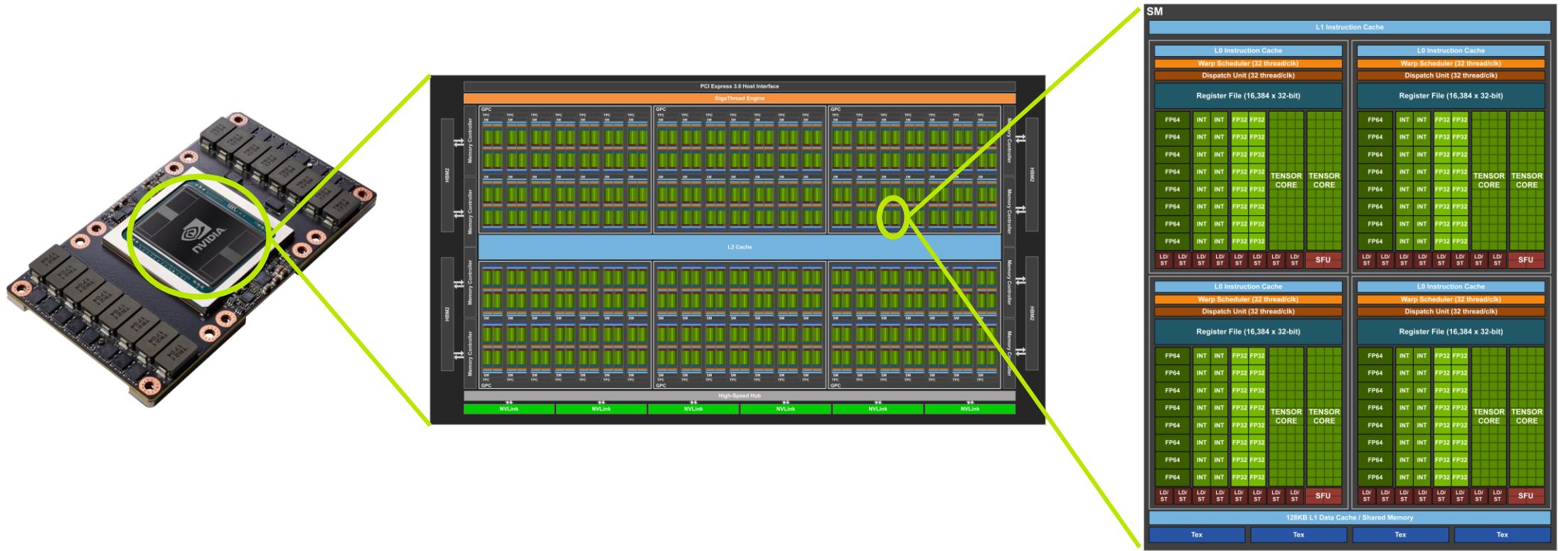
1. What are Tensor Cores?
2. Tensor Cores for AI
3. Tensor Cores for HPC
4. Profiling Tensor Cores

# FIRST, WHAT IS PRECISION?

- Precision is a measure of numerical detail
- Floating Point (FP) is a representation of real numbers supporting the tradeoff of:
  - Precision (**significand**)
  - Range (**exponent**)
- Lower precision numbers have computational performance advantages



# WHAT ARE TENSOR / CUDA CORES?



Figures: <https://images.nvidia.com/content/volta-architecture/pdf/volta-architecture-whitepaper.pdf>

# VOLTA GV100 SM

|                          | GV100  |
|--------------------------|--------|
| FP32 units               | 64     |
| FP64 units               | 32     |
| INT32 units              | 64     |
| Tensor Cores             | 8      |
| Register File            | 256 KB |
| Unified L1/Shared memory | 128 KB |
| Active Threads           | 2048   |

CUDA CORES



# VOLTA TENSOR CORE

Half/Mixed Precision 4x4 Matrix Multiply-Accumulate

$$\mathbf{D} = \mathbf{AB} + \mathbf{C}$$

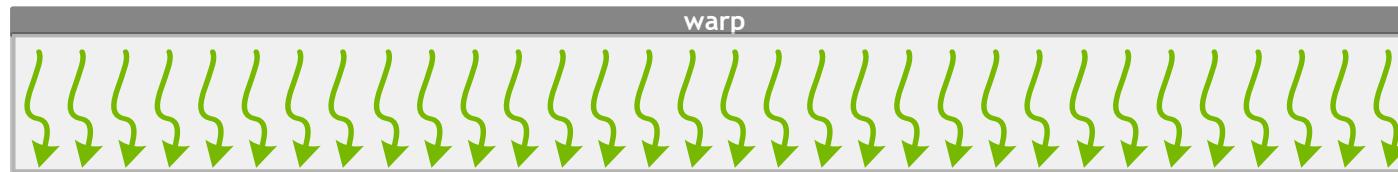
$$\mathbf{D} = \left( \begin{array}{cccc} \mathbf{A}_{0,0} & \mathbf{A}_{0,1} & \mathbf{A}_{0,2} & \mathbf{A}_{0,3} \\ \mathbf{A}_{1,0} & \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & \mathbf{A}_{1,3} \\ \mathbf{A}_{2,0} & \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} \\ \mathbf{A}_{3,0} & \mathbf{A}_{3,1} & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} \end{array} \right) \left( \begin{array}{cccc} \mathbf{B}_{0,0} & \mathbf{B}_{0,1} & \mathbf{B}_{0,2} & \mathbf{B}_{0,3} \\ \mathbf{B}_{1,0} & \mathbf{B}_{1,1} & \mathbf{B}_{1,2} & \mathbf{B}_{1,3} \\ \mathbf{B}_{2,0} & \mathbf{B}_{2,1} & \mathbf{B}_{2,2} & \mathbf{B}_{2,3} \\ \mathbf{B}_{3,0} & \mathbf{B}_{3,1} & \mathbf{B}_{3,2} & \mathbf{B}_{3,3} \end{array} \right) + \left( \begin{array}{cccc} \mathbf{C}_{0,0} & \mathbf{C}_{0,1} & \mathbf{C}_{0,2} & \mathbf{C}_{0,3} \\ \mathbf{C}_{1,0} & \mathbf{C}_{1,1} & \mathbf{C}_{1,2} & \mathbf{C}_{1,3} \\ \mathbf{C}_{2,0} & \mathbf{C}_{2,1} & \mathbf{C}_{2,2} & \mathbf{C}_{2,3} \\ \mathbf{C}_{3,0} & \mathbf{C}_{3,1} & \mathbf{C}_{3,2} & \mathbf{C}_{3,3} \end{array} \right)$$

FP16 or FP32                    FP16                    FP16 or FP32

Turing Tensor Cores: support for int8, int4

# VOLTA TENSOR CORE

Full Warp 16x16 Matrix Math

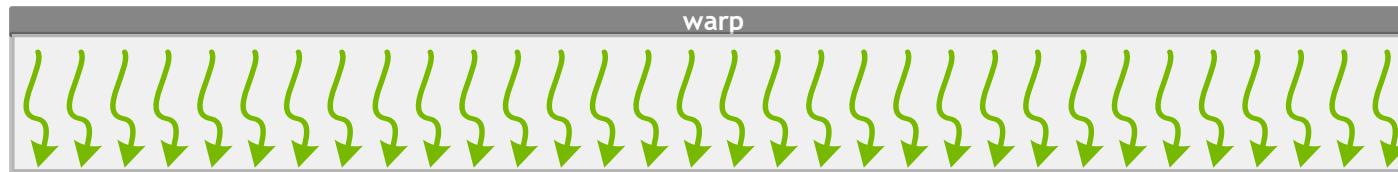


Warp-synchronizing operation for cooperative matrix math

$$\begin{pmatrix} \text{teal} \\ \text{purple} \end{pmatrix} \begin{pmatrix} \text{teal} \\ \text{purple} \end{pmatrix} + \begin{pmatrix} \text{green} \end{pmatrix}$$

$$\begin{pmatrix} \text{teal} \\ \text{purple} \end{pmatrix} \begin{pmatrix} \text{teal} \\ \text{purple} \end{pmatrix} + \begin{pmatrix} \text{green} \end{pmatrix}$$

Aggregate Matrix Multiply and Accumulate for 16x16 matrices



Result distributed across warp

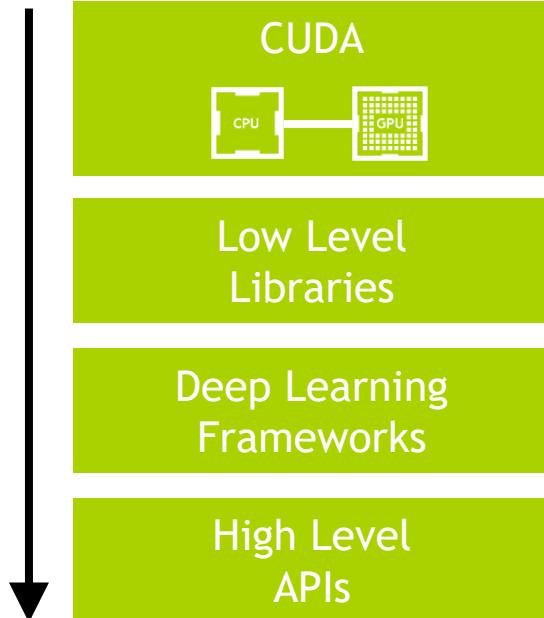
# OUTLINE

1. What are Tensor Cores?
2. **Tensor Cores for AI**
3. Tensor Cores for HPC
4. Profiling Tensor Cores

# TENSOR CORES FOR AI

- Simple trick for **2x to 5x faster deep learning training**
  - Accomplished in **few lines of code**
  - Models can use **same hyperparameters**
  - Models converge to **same accuracy**
- Half the memory traffic and storage **enabling larger batch sizes**
- AI community is trending towards **low precision as common practice**

# HOW TO USE TENSOR CORES



- Exposed as instructions in CUDA under WMMA API (**Warp Matrix Multiply Accumulate**)
- Used by cuDNN, cuBLAS, CUTLASS to accelerate matrix multiplications and convolution
- Tensor Core kernels used implicitly on FP16 ops from DL frameworks PyTorch / TensorFlow / etc...
- High-level tools (e.g. PyTorch Apex) convert everything automatically and safely

# MIXED PRECISION TRAINING

1. Model conversion
2. Master weight copy
3. Loss scaling

# 1. MODEL CONVERSION

- Make simple type updates to each layer:
  - Use FP16 values for the weights and inputs

```
# PyTorch
layer = torch.nn.Linear(in_dim, out_dim).half()

# TensorFlow
layer = tf.layers.dense(tf.cast(inputs, tf.float16), out_dim)
```

## 2. MASTER WEIGHTS

- FP16 alone is sufficient for some networks but not others; keep FP32 copy of weights

```
param = torch.cuda.FloatTensor([1.0])  
print(param + 0.0001)
```



1.0001

```
param = torch.cuda.HalfTensor([1.0])  
print(param + 0.0001)
```

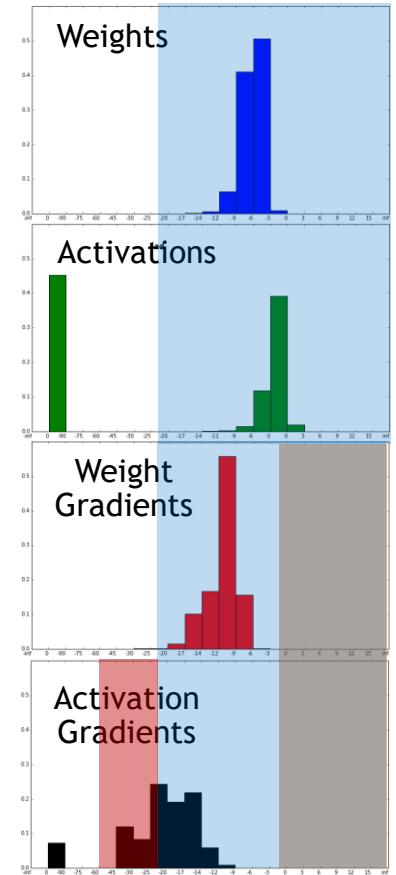


1

When *update/param < 2<sup>-11</sup>*, updates have no effect.

# 3. LOSS SCALING

- Range representable in FP16: ~40 powers of 2
- Gradients are small:
  - Some lost to zero
  - While ~15 powers of 2 remain unused
- Loss scaling:
  - Multiply loss by a constant  $S$
  - All gradients scaled up by  $S$  (chain rule)
  - Unscale weight gradient (in FP32) before weight update



# MIXED PRECISION TRAINING

1. Model conversion
2. Master weight copy
3. Loss scaling



Automated Mixed Precision (AMP)  
(e.g. PyTorch Apex)

# PYTORCH APEX AMP 1.0

```
N, D_in, D_out = 64, 1024, 512
x = Variable(torch.randn(N, D_in )).cuda()
y = Variable(torch.randn(N, D_out)).cuda()

model = torch.nn.Linear(D_in, D_out).cuda()
optimizer = torch.optim.SGD(model.parameters(), lr=1e-3)

for t in range(500):
    y_pred = model(x)
    loss = torch.nn.functional.mse_loss(y_pred, y)

    optimizer.zero_grad()

    loss.backward()
    optimizer.step()
```

# PYTORCH APEX AMP 1.0

```
N, D_in, D_out = 64, 1024, 512
x = Variable(torch.randn(N, D_in)).cuda()
y = Variable(torch.randn(N, D_out)).cuda()

model = torch.nn.Linear(D_in, D_out).cuda()
optimizer = torch.optim.SGD(model.parameters(), lr=1e-3)
model, optimizer = amp.initialize(model, optimizer, opt_level="O1")

for t in range(500):
    y_pred = model(x)
    loss = torch.nn.functional.mse_loss(y_pred, y)

    optimizer.zero_grad()
    with amp.scale_loss(loss, optimizer) as scaled_loss:
        scaled_loss.backward()
    optimizer.step()
```

# MIXED PRECISION SPEEDUPS

Not Limited to Image Classification

| Model                                   | FP32 -> FP16 Speedup | Comments                               |
|---|----------------------|--|
| GNMT<br>(Translation)                   | 2.3x                 | Iso-batch size                         |
| FairSeq Transformer<br>(Translation)    | 2.9x<br>4.9x         | Iso-batch size<br>2x lr + larger batch |
| ConvSeq2Seq<br>(Translation)            | 2.5x                 | 2x batch size                          |
| Deep Speech 2<br>(Speech recognition)   | 4.5x                 | Larger batch                           |
| wav2letter<br>(Speech recognition)      | 3.0x                 | 2x batch size                          |
| Nvidia Sentiment<br>(Language modeling) | 4.0x                 | Larger batch                           |

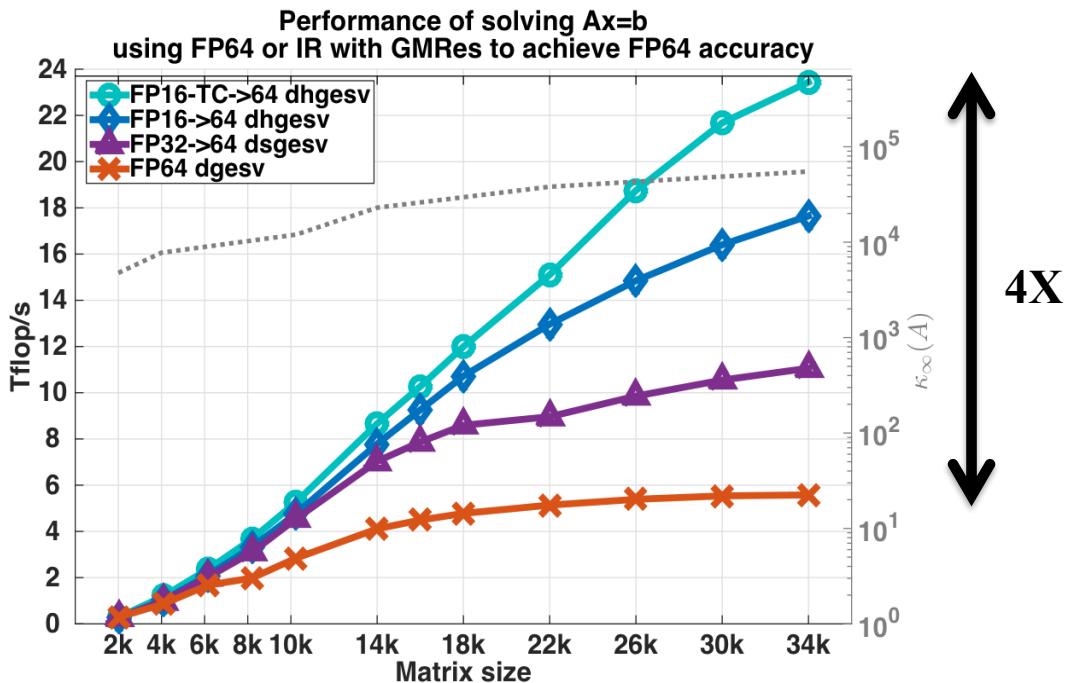
# OUTLINE

1. What are Tensor Cores?
2. Tensor Cores for AI
3. **Tensor Cores for HPC**
4. Profiling Tensor Cores

# TENSOR CORES FOR HPC

- Mixed precision algorithms are **increasingly popular**
  - It is common to combine **double + single precision**, or **floating point + integer**
- Similar to AI:
  - Use low precision to **reduce memory traffic and storage**
  - Use Tensor Core instructions for **large speedups**

# LINEAR ALGEBRA



- Researchers from ICL/UTK
- Accelerated FP64 LU factorization 4x using Tensor Cores in MAGMA
- Compute initial solution in FP16, then iteratively refine solution
- Achieved FP64 TFLOPS: 5.8
- Achieved FP16->FP64 TFLOPS: 24

Data courtesy of: Azzam Haidar, Stan. Tomov & Jack Dongarra, Innovative Computing Laboratory, University of Tennessee

"Harnessing GPU Tensor Cores for Fast FP16 Arithmetic to Speed up Mixed-Precision Iterative Refinement Solvers", A. Haidar, S. Tomov, J. Dongarra, N. Higham SC'18  
GTC 2018 Poster P8237: Harnessing GPU's Tensor Cores Fast FP16 Arithmetic to Speedup Mixed-Precision Iterative Solves

# EARTHQUAKE SIMULATION

- Researchers from University of Tokyo, Oak Ridge National Laboratory (ORNL), and the Swiss National Supercomputing Centre
- Solver called MOTHRA achieved 25x compared to standard solver
- Used AI to identify where to apply low or high precision in solver
- Used a combination FP64, FP32, FP21 and FP16 to further reduce computational and communication costs

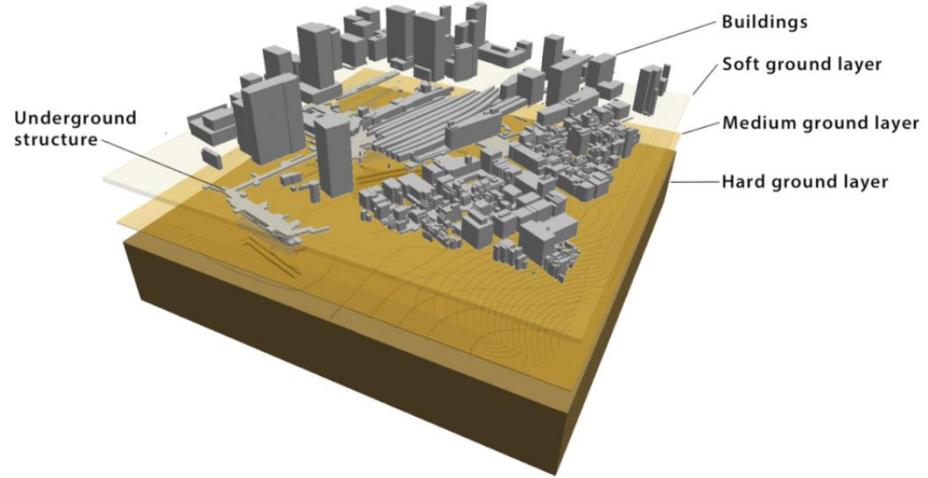


Figure 3. Earthquake simulation model used by the Supercomputing 2018 (SC'18) Gordon Bell finalist paper for Tokyo including layers of soil, underground and above ground structures, and the coupling between them.

# OUTLINE

1. What are Tensor Cores?
2. Tensor Cores for AI
3. Tensor Cores for HPC
4. Profiling Tensor Cores

# NSIGHT DEVELOPER TOOLS

# NSIGHT PRODUCT FAMILY

## Standalone Performance Tools

**Nsight Systems** - System-wide application algorithm tuning

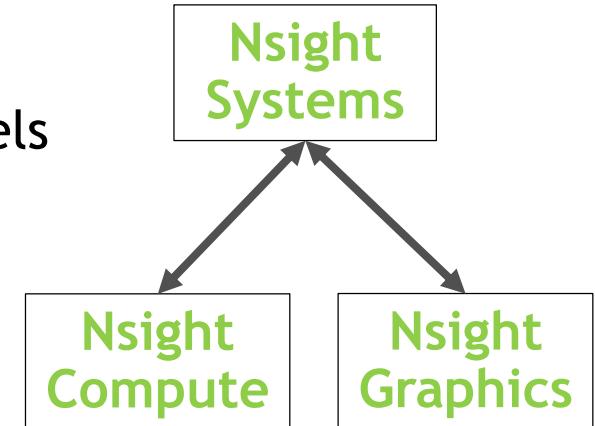
**Nsight Compute** - Debug CUDA API and optimize CUDA kernels

**Nsight Graphics** - Debug/optimize specific graphics apps

## IDE Plugins

**Nsight Eclipse Edition/Visual Studio** - editor, debugger, some perf analysis

## Workflow





# NSIGHT SYSTEMS

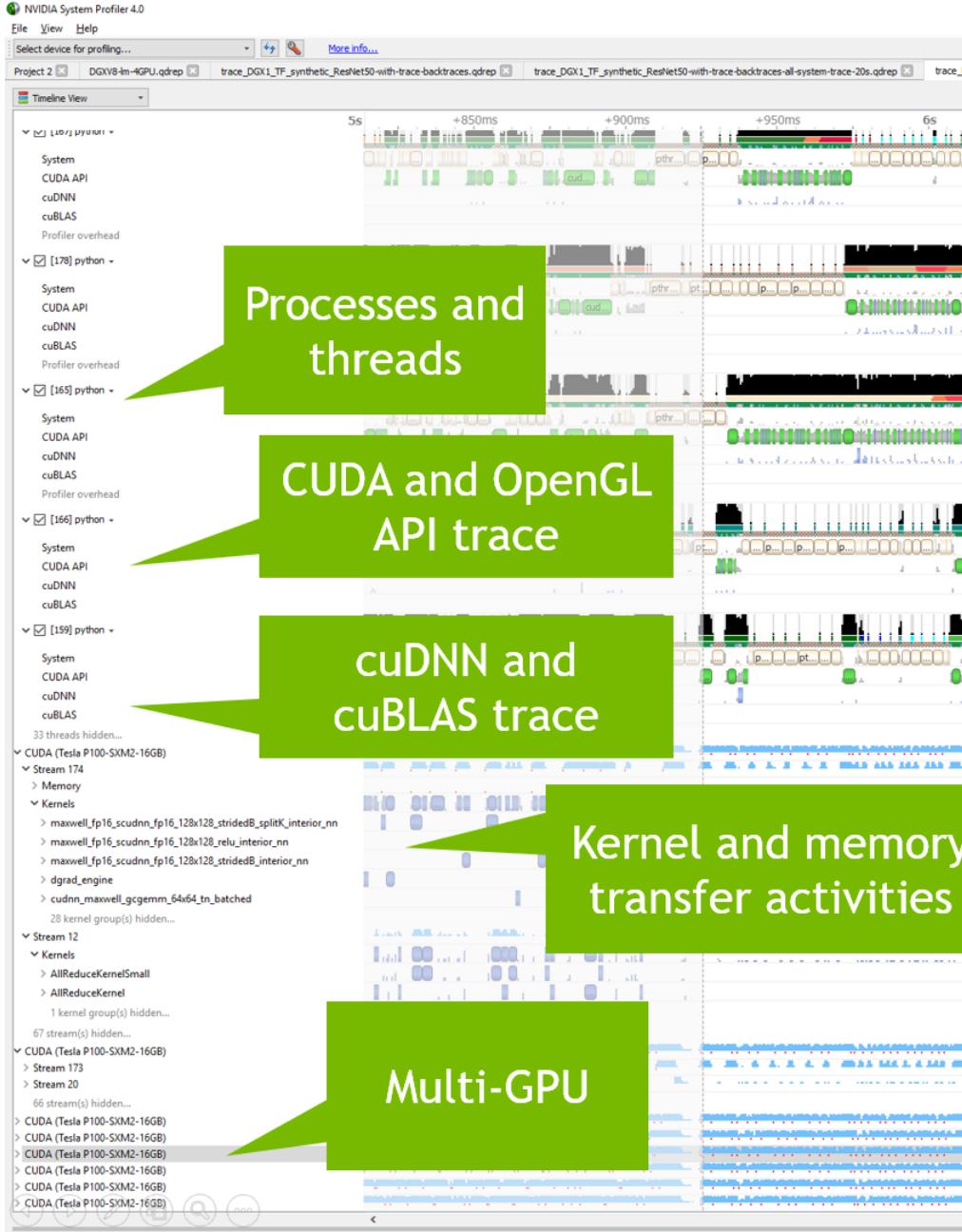
## Next-Gen System Profiling Tool

System-wide application algorithm tuning  
Multi-process tree support

Locate optimization opportunities  
Visualize millions of events on a fast GUI timeline  
Or gaps of unused CPU and GPU time

Balance your workload across multiple CPUs and GPUs  
CPU algorithms, utilization, and thread state  
GPU streams, kernels, memory transfers, etc

Multi-platform: Linux & Windows, x86-64 & Tegra,  
MacOSX (host only)





# NSIGHT COMPUTE

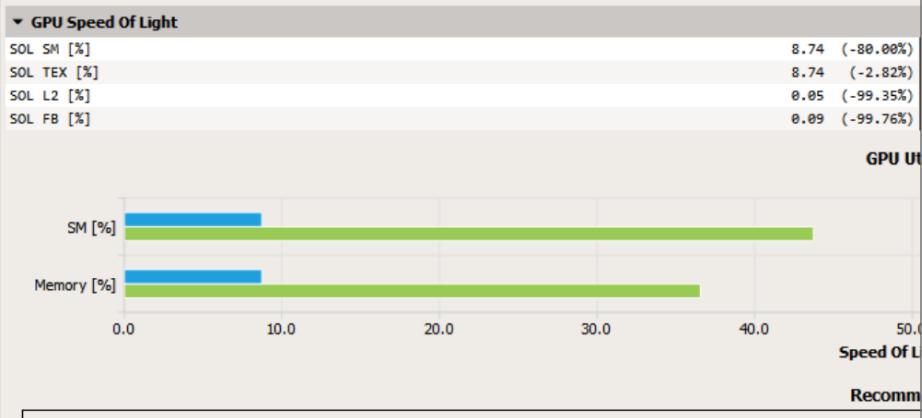
## Next-Gen Kernel Profiling Tool

### Key Features:

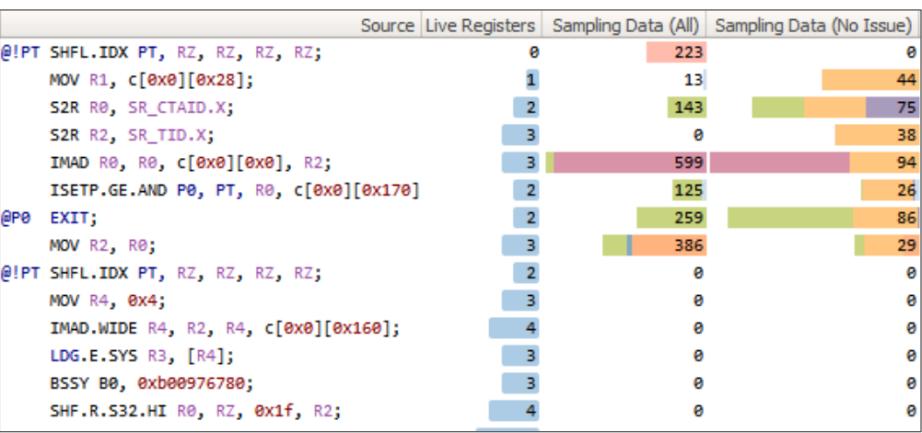
- Interactive CUDA API debugging and kernel profiling
- Fast Data Collection
- Improved Workflow (differencing results)
- Fully Customizable (programmable UI/Rules)
- Command Line, Standalone, IDE Integration

OS: Linux, Windows, ARM, MacOSX (host only)

GPUs: Pascal (GP10x), Volta, Turing



|   |                           |                           |                           |                           |
|---|---------------------------|---------------------------|---------------------------|---------------------------|
| inst_executed [inst]                              | 16,528.00; 16,528.00; ... | 13,476.00; 13,476.00; ... | 14.33                     | n/a                       |
| l1tex_solid_pct [%]                               |                           |                           | 128.00                    | 128.00                    |
| launch_block_size                                 |                           |                           | 47,611,587,968.00         | 12,273,728.00             |
| launch_function_pcs                               |                           |                           | 4,132.00                  | 3,369.00                  |
| launch_grid_size                                  |                           |                           | 32.00                     | 32.00                     |
| launch_occupancy_limit_blocks [block]             |                           |                           | 21.00                     | 21.00                     |
| launch_occupancy_limit_registers [register]       |                           |                           | 384.00                    | 384.00                    |
| launch_occupancy_limit_shared_mem [bytes]         |                           |                           | 16.00                     | 16.00                     |
| launch_occupancy_limit_warp [warps]               |                           |                           | 3,638.00                  | 3,638.00                  |
| launch_occupancy_per_block_size                   |                           |                           | 5,792.00                  | 5,792.00                  |
| launch_occupancy_per_register_count               |                           |                           | 2,260.00                  | 2,260.00                  |
| launch_occupancy_per_shared_mem_size              |                           |                           | 17.00                     | 17.00                     |
| launch_registers_per_thread [register/thread]     |                           |                           | 49,152.00                 | 49,152.00                 |
| launch_shared_mem_config_size [bytes]             |                           |                           | 0.00                      | 0.00                      |
| launch_shared_mem_per_block_dynamic [bytes/block] |                           |                           | 20.00                     | 20.00                     |
| launch_shared_mem_per_block_static [bytes/block]  |                           |                           | 528,896.00                | 431,232.00                |
| launch_thread_count [thread]                      |                           |                           | 3.23                      | 42.11                     |
| launch_waves_per_multiprocessor                   |                           |                           | 6.93                      | 7.18                      |
| ltc_solid_pct [%]                                 |                           |                           | 2.00; 32.00; 32.00; 32.00 | 2.00; 32.00; 32.00; 32.00 |
| memory_access_size_type [bytes]                   |                           |                           | 3.00; 1.00; 1.00; 1.00    | 3.00; 1.00; 1.00; 1.00    |



# USING NSIGHT SYSTEMS

# COLLECT A PROFILE WITH NSIGHT SYSTEMS

```
$ nsys profile /usr/bin/python train.py
```

Generated file: report.qdrep

Import for viewing into the Nsight Systems UI

The Nsight Systems UI can also be used for interactive system profiling

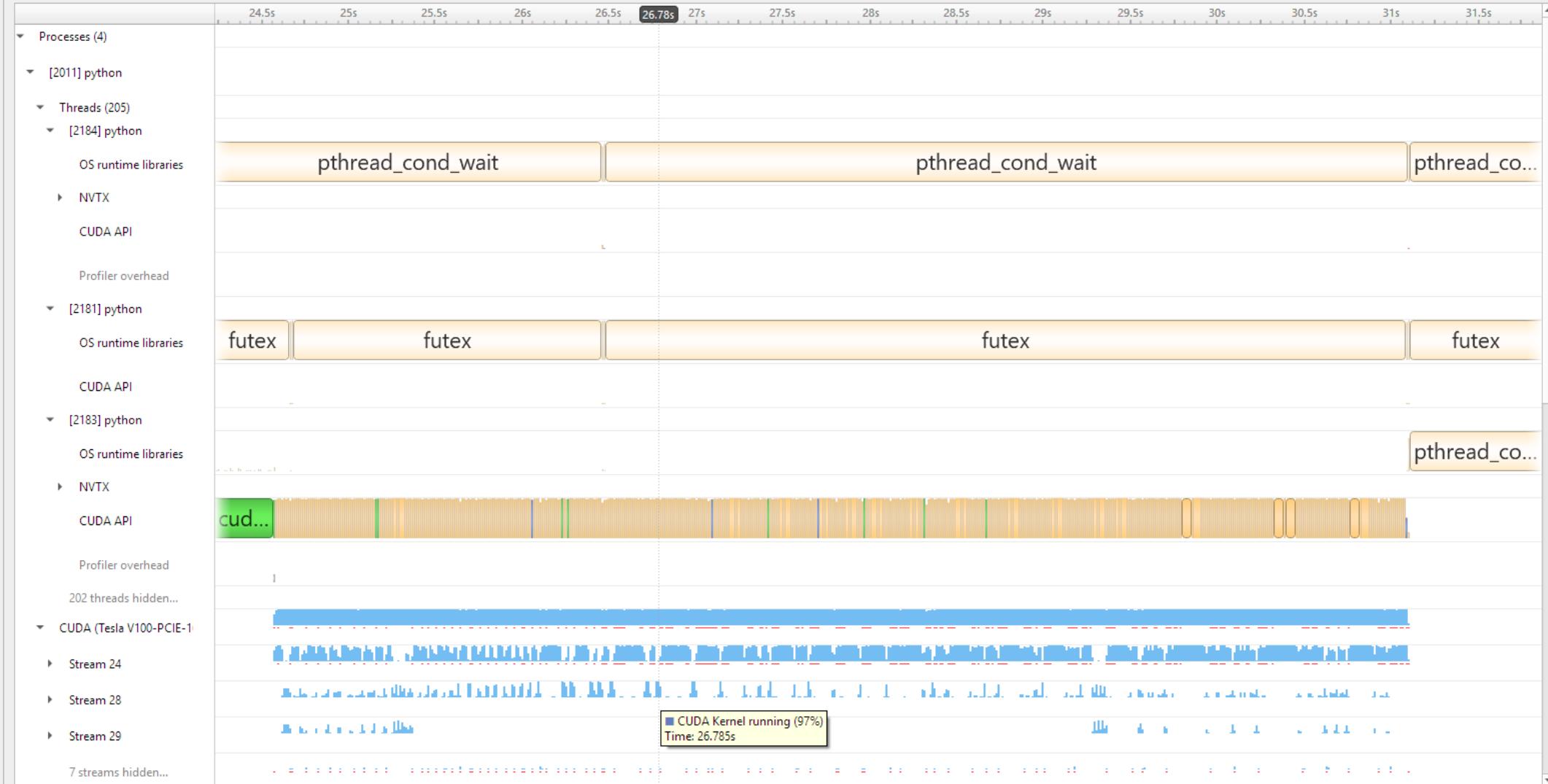
File View Tools Help

rn50.qdrep X

Timeline View

2x

2 errors, 19 warnings, 26 messages



# LOCATING TENSOR CORE KERNELS

On Volta V100, CUDA kernels using tensor cores contain the string “s884”

Examples:

volta\_fp16\_s884gemm\_fp16\_128x64\_ldg8\_f2f\_nn

volta\_fp16\_s884cudnn\_fp16\_256x64\_ldg8\_relu\_f2f\_exp\_interior\_nhwc2nchw\_tn\_v1

These are kernels with HMMA (half-precision matrix multiply and accumulate) machine instructions

File View Tools Help

rn50.qdrep X

Timeline View

1x

2 errors, 19 warnings, 26 messages

- ▼ [2184] python
  - OS runtime libraries

- ▶ NVTX
  - CUDA API
  - Profiler overhead

- ▼ [2181] python
  - OS runtime libraries
  - CUDA API

- ▼ [2183] python
  - OS runtime libraries
  - NVTX
    - CUDA API
    - Profiler overhead

202 threads hidden...

▼ CUDA (Tesla V100-PCIE-16GB, 0000:04:00.0)

▼ Stream 24

▼ Kernels

- ▶ dgrad\_engine
- ▶ wgrad\_alg0\_engine
- ▶ volta\_cgemm\_32x32\_tn

- ▶ volta\_fp16\_scudnn\_fp16\_128x128\_stridedB\_splitK\_interior\_r
- ▶ nchwToNhwcKernel

- ▶ EigenMetaKernel

- ▶ winograd3x3Kernel

- ▶ explicit\_convolve\_sgemm

- ▶ implicit\_convolve\_sgemm

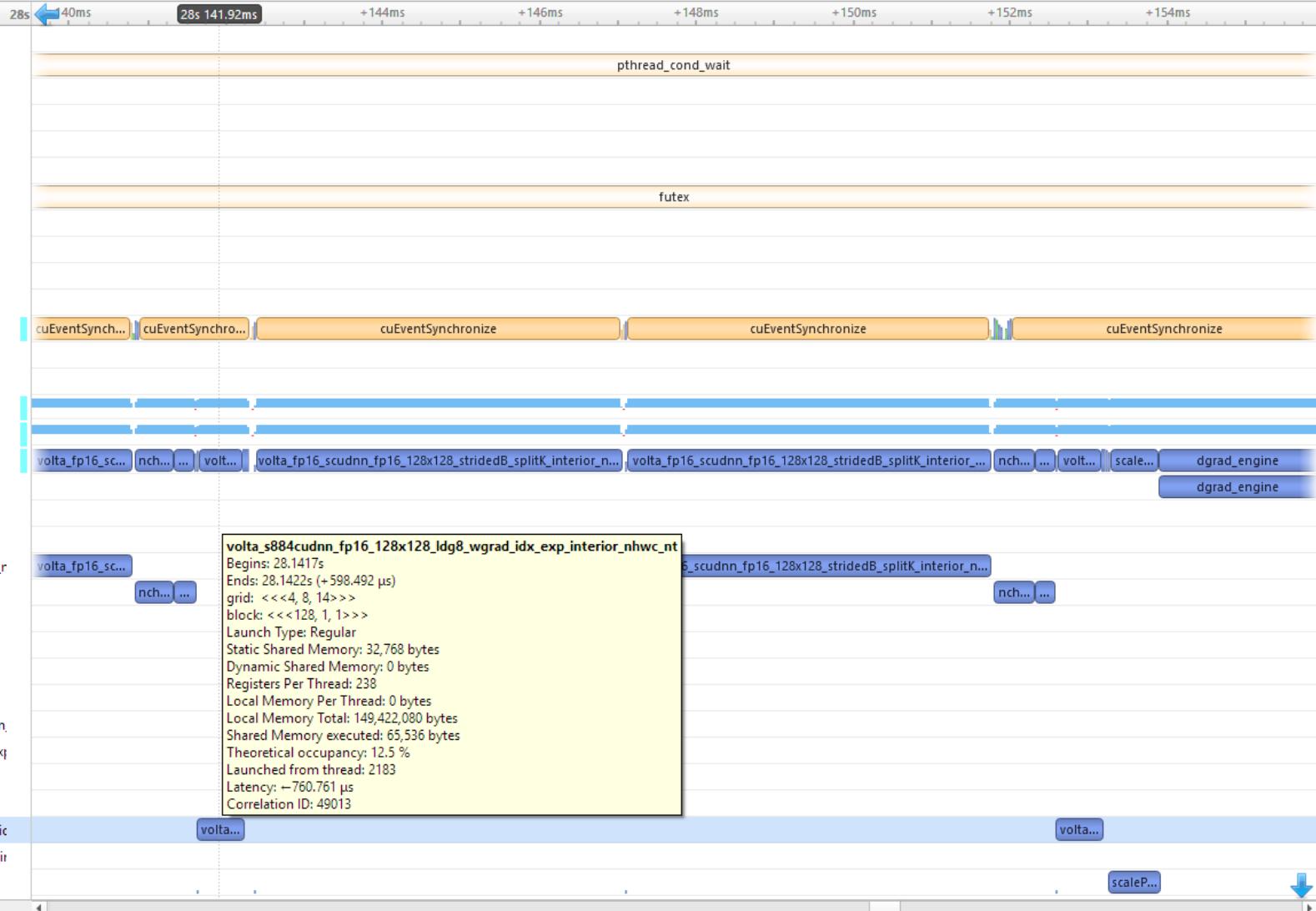
- ▶ volta\_fp16\_scudnn\_fp16\_128x128\_stridedB\_splitK\_small\_nn
- ▶ volta\_s884cudnn\_fp16\_64x64\_sliced1x4\_ldg8\_wgrad\_idx\_exq

- ▶ transpose\_readWrite\_alignment\_kernel

- ▶ ReluGradHalfKernel

- ▶ volta\_s884cudnn\_fp16\_128x128\_ldg8\_wgrad\_idx\_exp\_interic
- ▶ volta\_fp16\_s884cudnn\_fp16\_128x128\_ldg8\_dgrad\_f2f\_exp\_ir

- ▶ scalePackedTensor\_kernel



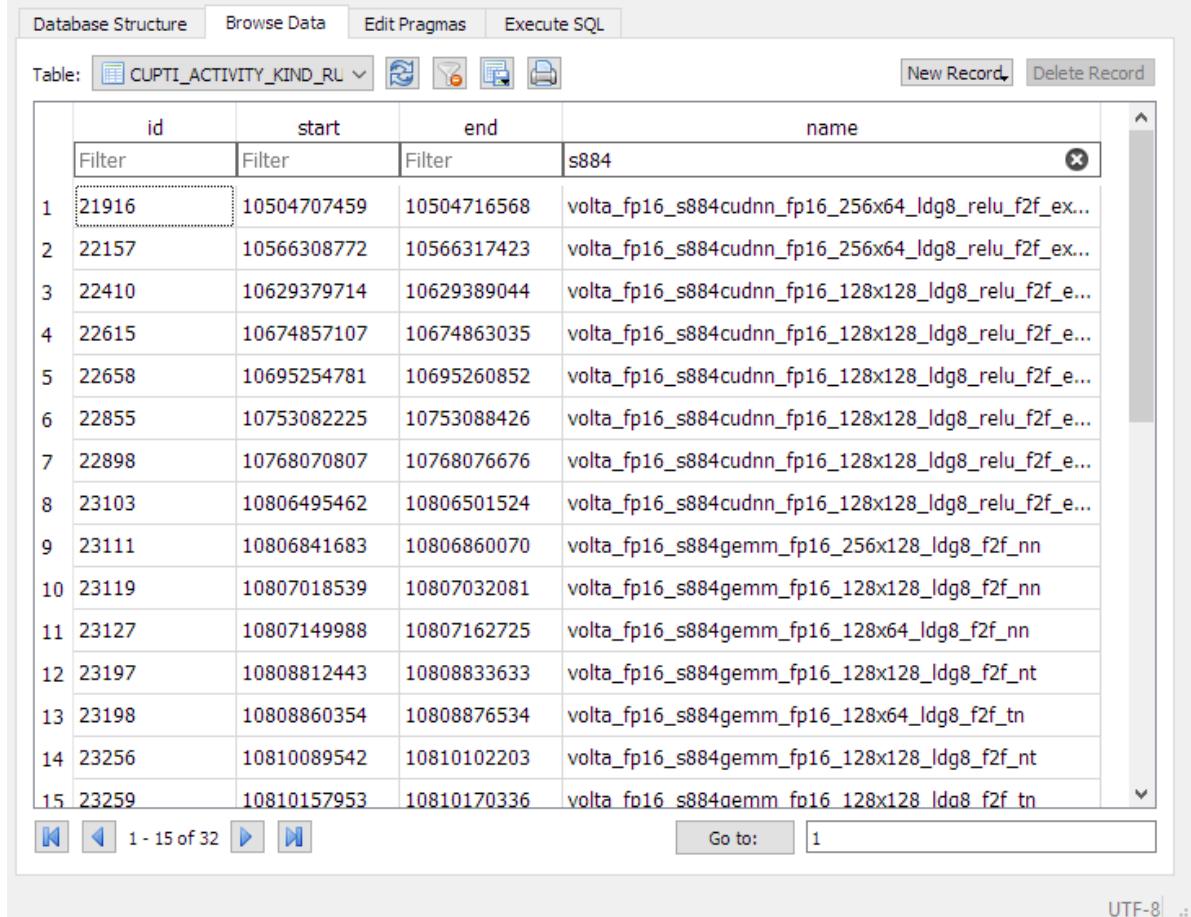
**COMING SOON: SQLITE DATABASE EXPORT**

# USE NSYS-EXPORTER TO CREATE SQLITE DB

```
nsys-exporter -s report.qdrep
```

Generated DB: report.sqlite

Interact with this like any SQLite database



The screenshot shows a SQLite database browser window with the following details:

- Toolbar buttons: Database Structure, Browse Data, Edit Pragmas, Execute SQL.
- Table selection: CUPTI\_ACTIVITY\_KIND\_RL.
- Action buttons: New Record, Delete Record.
- Table structure:

|    | id     | start       | end         | name   |
|----|--------|-------------|-------------|--|
|    | Filter | Filter      | Filter      | s884   |
| 1  | 21916  | 10504707459 | 10504716568 | volta_fp16_s884cudnn_fp16_256x64_ldg8_relu_f2f_ex... |
| 2  | 22157  | 10566308772 | 10566317423 | volta_fp16_s884cudnn_fp16_256x64_ldg8_relu_f2f_ex... |
| 3  | 22410  | 10629379714 | 10629389044 | volta_fp16_s884cudnn_fp16_128x128_ldg8_relu_f2f_e... |
| 4  | 22615  | 10674857107 | 10674863035 | volta_fp16_s884cudnn_fp16_128x128_ldg8_relu_f2f_e... |
| 5  | 22658  | 10695254781 | 10695260852 | volta_fp16_s884cudnn_fp16_128x128_ldg8_relu_f2f_e... |
| 6  | 22855  | 10753082225 | 10753088426 | volta_fp16_s884cudnn_fp16_128x128_ldg8_relu_f2f_e... |
| 7  | 22898  | 10768070807 | 10768076676 | volta_fp16_s884cudnn_fp16_128x128_ldg8_relu_f2f_e... |
| 8  | 23103  | 10806495462 | 10806501524 | volta_fp16_s884cudnn_fp16_128x128_ldg8_relu_f2f_e... |
| 9  | 23111  | 10806841683 | 10806860070 | volta_fp16_s884gemm_fp16_256x128_ldg8_f2f_nn         |
| 10 | 23119  | 10807018539 | 10807032081 | volta_fp16_s884gemm_fp16_128x128_ldg8_f2f_nn         |
| 11 | 23127  | 10807149988 | 10807162725 | volta_fp16_s884gemm_fp16_128x64_ldg8_f2f_nn          |
| 12 | 23197  | 10808812443 | 10808833633 | volta_fp16_s884gemm_fp16_128x128_ldg8_f2f_nt         |
| 13 | 23198  | 10808860354 | 10808876534 | volta_fp16_s884gemm_fp16_128x64_ldg8_f2f_nt          |
| 14 | 23256  | 10810089542 | 10810102203 | volta_fp16_s884gemm_fp16_128x128_ldg8_f2f_nt         |
| 15 | 23259  | 10810157953 | 10810170336 | volta_fp16_s884gemm_fp16_128x128_ldg8_f2f_nt         |
- Pagination: 1 - 15 of 32.
- Search: Go to: 1.
- Encoding: UTF-8.

# ASSOCIATE KERNEL NAMES WITH EVENTS

```
ALTER TABLE CUPTI_ACTIVITY_KIND_RUNTIME ADD COLUMN name TEXT;  
  
UPDATE CUPTI_ACTIVITY_KIND_RUNTIME SET name =  
    (SELECT value FROM StringIds  
     JOIN CUPTI_ACTIVITY_KIND_KERNEL AS cuda_gpu  
       ON cuda_gpu.demangledName = StringIds.id  
      AND CUPTI_ACTIVITY_KIND_RUNTIME.correlationId =  
        cuda_gpu.correlationId);
```

# LOCATE KERNELS USING TENSOR CORES

SELECT

\*

FROM

CUPTI\_ACTIVITY\_KIND\_RUNTIME as cupti

WHERE

cupti.name LIKE '%s884%'

# USING NSIGHT COMPUTE

# KERNEL PROFILES WITH NSIGHT COMPUTE

```
$ nv-nsight-cu-cli /usr/bin/python train.py
```

This is expensive for a typical DL training session because it will collect metrics for every kernel; consider profiling fewer kernels.

For example, to profile \*s884\* kernels on all streams, but only on the fifth invocation:

```
$ nv-nsight-cu-cli --kernel-id ::s884:5 /usr/bin/python train.py
```

The Nsight Systems UI can also be used for interactive kernel profiling

# INTERLUDE: FINDING TENSOR CORE METRICS

Isolating which GPU metrics measure tensor cores:

```
$ nv-nsight-cu-cli --devices 0 --query-metrics | grep -i tensor  
...  
smsp_pipe_tensor_cycles_active.avg.pct_of_peak_sustained_active  
sm_pipe_tensor_op_hmma_cycles_active.avg  
sm_inst_executed_pipe_tensor.avg.per_second  
...
```

# TENSOR CORE PERFORMANCE

Now add in what we know about which metrics to looks for:

```
$ nv-nsight-cu-cli --kernel-id ::s884:5 --metrics  
smsp__pipe_tensor_cycles_active.avg.pct_of_peak_sustained_active  
/usr/bin/python train.py
```

```
volta_s884cudnn_fp16_128x128_ldg8_wgrad_idx_exp_interior_nhwc_nt, 2019-Mar-17 02:54:29, Context 1, Stream 23
```

```
Section: Command line profiler metrics
```

---

|  |   |       |
|--|---|-------|
| smsp__pipe_tensor_cycles_active.avg.pct_of_peak_sustained_active | % | 79.03 |
|--|---|-------|

---

# NSIGHT COMPUTE UI

Current 47828 - volta\_fp16\_s884cudnn\_fp16\_256x128\_lg8\_relu\_filter1x1\_stg8\_interior\_nchw\_nn\_v1 (802816, 2, 1) Time: n/a Cycles: n/a Regs: 226 GPU: Tesla V100-PCIE-16GB SM Frequency: n/a CC: 7.0 +

GPU Speed Of Light

| SOL SM [%]     | Duration [msecond]       | 1.20             |
|----------------|--------------------------|------------------|
| SOL Memory [%] | Elapsed Cycles [cycle]   | 1,504,418        |
| SOL TEX [%]    | SM Active Cycles [cycle] | 1,459,787.71     |
| SOL L2 [%]     | SM Frequency [cycle]     | 1,248,295,701.36 |
| SOL FB [%]     | Memory Frequency [cycle] | 881,625,488.68   |

GPU Utilization

| Category   | Utilization [%] |
|------------|-----------------|
| SM [%]     | ~28             |
| Memory [%] | ~57             |

Recommendations

Bottleneck High-level bottleneck detection Apply

Compute Workload Analysis

|                              |      |                      |       |
|------------------------------|------|----------------------|-------|
| Executed Ipc Elapsed [cycle] | 1.11 | SM Busy [%]          | 28.44 |
| Executed Ipc Active [cycle]  | 1.14 | Issue Slots Busy [%] | 28.44 |
| Issued Ipc Active [cycle]    | 1.14 | Issue Slots Max [%]  | 29.02 |

Memory Workload Analysis

|                           |        |                    |       |
|---------------------------|--------|--------------------|-------|
| Memory Throughput [Gbyte] | 477.73 | Mem Busy [%]       | 52.23 |
| L1 Hit Rate [%]           | 13.35  | Max Bandwidth [%]  | 56.82 |
| L2 Hit Rate [%]           | 69.23  | Mem Pipes Busy [%] | 15.33 |

Scheduler Statistics

|                                     |      |   |       |
|-------------------------------------|------|---|-------|
| Active Warps Per Scheduler [warp]   | 2.01 | Instructions Per Active Issue Slot [inst/issue] | 1     |
| Eligible Warps Per Scheduler [warp] | 0.40 | No Eligible [%]                                 | 71.48 |
| Issued Warp Per Scheduler [cycle]   | 0.29 | One or More Eligible [%]                        | 28.52 |

# SUMMARY

# NVIDIA TOOLS FOR TENSOR CORE PROFILING

System, application, and kernel level profiling solutions

- **Nsight Systems:** high-level application view; locate kernels that used tensor cores
- **Nsight Compute:** drill down into specific kernels for detailed performance analysis
- Starting with version 19.03, [NVIDIA GPU Cloud](#) (NGC) optimized deep learning containers package Nsight Systems and Nsight Compute. Download and try it now!
- Coming soon: DLProf tool for connecting tensor core usage with popular DL frameworks. More information: [S9339 - Profiling Deep Learning Networks @ GTC 2019](#)

# DEVELOPER TOOLS AT GTC19

## Talks:

- S9751: Accelerate Your CUDA Development with Latest Debugging and Code Analysis Developer Tools, **Tue @9am**
- S9866 - Optimizing Facebook AI Workloads for NVIDIA GPUs, **Tue @9am**
- S9345: CUDA Kernel Profiling using NVIDIA Nsight Compute, **Tue @1pm**
- S9661: Nsight Graphics - DXR/Vulkan Profiling/Vulkan Raytracing, **Wed @10am**
- S9503: Using Nsight Tools to Optimize the NAMD Molecular Dynamics Simulation Program, **Wed @1pm**

## Hands-on labs:

- L9102: Jetson Developer Tools Training Lab, **Mon @9am, 11:30am**
- L9124: Debugging and optimizing CUDA applications with Nsight products on Linux training lab, **Tue @8am, 10am**

## Connect with the Experts (where DevTools will be available):

- CE9123: CUDA & Graphics Developer Tools, **Tue @2pm, Wed @3pm**
- CE9137: Jetson Embedded Platform, **Tue @12pm, 5pm, Wed @1pm, 4pm, Thu @12pm**

## Podium: Demos of DevTools products on Linux, DRIVE AGX & Jetson AGX at the showfloor

- Tue @12pm - 7pm**
- Wed @12pm - 7pm**
- Thu @11am - 2pm**

PRESENTED BY  
 NVIDIA

# GPU TECHNOLOGY CONFERENCE

SAN JOSE | MAY 17-21, 2019

<https://www.nvidia.com/en-us/gtc/>

#GTC19